**Title: A Machine Learning Approach for Optimal Crop Recommendation Based on Environmental and Soil Parameters**

*Sheheryar Yousaf /Matricola: B33000111*

**Abstract:**

Precision agriculture aims to optimize farming practices by leveraging data-driven insights. Selecting the appropriate crop for specific environmental and soil conditions is crucial for maximizing yield and ensuring agricultural sustainability. This study presents a machine learning-based approach for recommending optimal crops based on key parameters: Nitrogen (N), Phosphorus (P), Potassium (K) content in the soil, temperature, humidity, pH level, and rainfall. Utilizing a publicly available dataset containing 2200 instances across 22 different crop types, we explored various classification models, including Random Forest, Logistic Regression, Support Vector Machines (SVM), and K-Nearest Neighbors (KNN). Data preprocessing involved feature scaling and appropriate encoding. Model performance was evaluated using accuracy, precision, recall, and F1-score. The Random Forest classifier, after hyperparameter optimization using RandomizedSearchCV, demonstrated superior performance, achieving near-perfect accuracy (approximately 99.5% on cross-validation and 100% on the test set) in recommending suitable crops. This research highlights the potential of machine learning models to serve as effective decision-support tools for farmers, enhancing agricultural productivity and resource management.

**Keywords:** Crop Recommendation, Machine Learning, Precision Agriculture, Random Forest, Environmental Parameters, Soil Nutrients.

## 1. Introduction

Agriculture is the backbone of global food security and economic stability. However, traditional farming practices often face challenges related to suboptimal crop selection, leading to reduced yields and inefficient resource utilization [1]. The decision of which crop to cultivate is complex, influenced by a multitude of interacting environmental factors such as climate conditions (temperature, humidity, rainfall) and soil properties (nutrient content like N, P, K, and pH levels) [2]. In an era of increasing climate variability and the need for sustainable agriculture, data-driven decision-making tools are becoming indispensable.

Precision agriculture leverages technologies like machine learning (ML) to analyze complex datasets and provide actionable insights for optimizing farming operations [3]. ML algorithms can learn patterns from historical data, relating crop performance to specific environmental and soil characteristics. This capability can be harnessed to develop recommendation systems that guide farmers in selecting the most suitable crops for their land, thereby potentially increasing productivity, reducing waste, and promoting sustainable land use.

This study aims to develop and evaluate a machine learning model for recommending optimal crops based on seven key environmental and soil parameters: Nitrogen, Phosphorus, Potassium, temperature, humidity, pH, and rainfall. We investigate the performance of several

common classification algorithms and identify the most effective model for this task, demonstrating its potential as a practical tool for agricultural decision support.

## 2. Materials and Methods

### 2.1. Dataset Description
The study utilized the "Crop_recommendation.csv" dataset, a publicly available collection of agricultural data. The dataset comprises 2200 instances, each representing a specific set of environmental and soil conditions, along with the corresponding recommended crop. It includes the following eight features:

- **N:** Ratio of Nitrogen content in soil (integer).
- **P:** Ratio of Phosphorus content in soil (integer).
- **K:** Ratio of Potassium content in soil (integer).
- **temperature:** Temperature in degrees Celsius (float).
- **humidity:** Relative humidity in percentage (float).
- **ph:** pH value of the soil (float).
- **rainfall:** Rainfall in mm (float).
- **label:** The recommended crop type (categorical, 22 unique crop types).

### 2.2. Data Preprocessing and Exploration
Initial data exploration was performed using Python libraries (Pandas, Matplotlib, Seaborn) within a Jupyter Notebook environment. This involved:

- **Data Loading and Inspection:** The dataset was loaded into a Pandas DataFrame. The shape of the data, data types of each column, and a summary of descriptive statistics (mean, std, min, max, quartiles) for numerical features were examined.
- **Target Variable Distribution:** The distribution of the target variable ('label') was analyzed to understand the frequency of each of the 22 crop types, confirming a balanced dataset with 100 instances per crop.
- **Correlation Analysis:** A correlation matrix and heatmap were generated for the numerical features to understand their interrelationships.
- **Missing Value Check:** The dataset was checked for missing values; none were found.

Data preprocessing steps included:

- **Feature Scaling:** Numerical features (N, P, K, temperature, humidity, ph, rainfall) were scaled using StandardScaler from Scikit-learn to ensure that all features contribute equally to model training, preventing features with larger magnitudes from dominating those with smaller magnitudes.
- **Categorical Data Handling:** The target variable 'label' was one-hot encoded to convert categorical crop names into a numerical format suitable for the machine learning algorithms. This resulted in 22 binary columns, each representing a crop type.

### 2.3. Feature Engineering (Attempted)
An attempt was made to improve model performance through feature engineering. New features

were created based on ratios, differences, and products of existing numerical features (e.g., 'N_P_ratio', 'temp_humidity_diff', 'N_K_product'). However, initial evaluations using a RandomForestClassifier indicated that these engineered features did not improve, and in some cases slightly degraded, model performance. Consequently, these engineered features were discarded, and the models were trained on the original scaled features.

### 2.4. Data Splitting

The preprocessed dataset was split into features (X) and the one-hot encoded target (y). The data was then divided into training, validation, and testing sets. An initial split allocated 70% of the data for training (X_train, y_train) and 30% for a temporary set (X_temp, y_temp). The temporary set was further split equally into validation (X_val, y_val) and testing sets (X_test, y_test), resulting in a 70% training, 15% validation, and 15% testing distribution. Stratification based on the original 'label' column was used during splitting to ensure that the proportion of crop types was maintained across all sets.

### 2.5. Model Selection and Training

Several classification algorithms were selected for evaluation:

- Random Forest Classifier
- Logistic Regression
- Support Vector Machine (SVC)
- K-Nearest Neighbors (KNN)

The target variable (one-hot encoded crop labels) was converted back to a 1D array of integer labels (0 to 21) using np.argmax for training and evaluation with Scikit-learn classifiers that expect 1D target arrays. All models were trained on the X_train and y_train_labels data.

### 2.6. Hyperparameter Optimization

The Random Forest classifier, which showed promising initial performance, was selected for hyperparameter optimization. RandomizedSearchCV from Scikit-learn was employed with 5-fold cross-validation. The parameter grid explored included:

- n_estimators: [50, 100, 200]
- max_depth: [None, 10, 20, 30]
- min_samples_split: [2, 5, 10]
- min_samples_leaf: [1, 2, 4]
- max_features: ['sqrt', 'log2', None]

The search was configured for 10 iterations, evaluating based on accuracy.

### 2.7. Model Evaluation

The performance of all trained models (both default and optimized) was evaluated on the validation set (X_val, y_val_labels). The final optimized Random Forest model was then evaluated on the unseen test set (X_test, y_test_labels). Evaluation metrics included:

- Accuracy

- Precision (weighted average)
- Recall (weighted average)
- F1-Score (weighted average)
- Confusion Matrix
- Classification Report

## 3. Results

### 3.1. Exploratory Data Analysis

The dataset was found to be well-balanced with 100 samples for each of the 22 crop types. No missing values were present. Descriptive statistics revealed varying ranges for different features, underscoring the need for feature scaling. Correlation analysis showed some expected relationships between environmental variables, but no prohibitively high multicollinearity that would necessitate feature removal.

### 3.2. Model Performance Comparison

Initial evaluation of the four selected models on the validation set (after training on X_train) yielded the following F1-scores (weighted average):

- Random Forest: 0.9909
- Logistic Regression: 0.9692
- Support Vector Machine (SVC): 0.9878
- K-Nearest Neighbors (KNN): 0.9755

The Random Forest classifier exhibited the highest F1-score and accuracy, making it the primary candidate for further optimization.

### 3.3. Hyperparameter Optimization Results

RandomizedSearchCV for the RandomForestClassifier identified the following best hyperparameters:

- n_estimators: 100
- min_samples_split: 10
- min_samples_leaf: 1
- max_features: 'log2'
- max_depth: 20

The best cross-validation accuracy achieved with these parameters was approximately 0.9948.

### 3.4. Optimized Model Evaluation on Test Set

The optimized RandomForestClassifier (best_rf_model) was evaluated on the unseen test set. The model achieved outstanding performance:

- Accuracy: 1.00 (near perfect, with minor misclassifications for 'chickpea' which was 0.94 due to 1 misclassification out of 15 test samples)
- Precision (weighted): 1.00

- Recall (weighted): 1.00
- F1-Score (weighted): 1.00

The classification report (Table 1 - *example, actual detailed table would be here*) showed perfect or near-perfect scores for all individual crop classes. The confusion matrix (Figure 1 - *example, actual figure would be here*) indicated a diagonal matrix with very few off-diagonal elements, visually confirming the high classification accuracy. For example, out of 330 test samples, there was only one instance of 'chickpea' being misclassified as 'kidneybeans'.

*(Ideally, include a table for the classification report and a figure for the confusion matrix here if this were a full paper submission)*

## 4. Discussion

The results of this study demonstrate the strong potential of machine learning, particularly the Random Forest algorithm, for providing accurate crop recommendations based on environmental and soil parameters. The optimized Random Forest model achieved near-perfect performance on the test set, suggesting its robustness and ability to generalize to unseen data.

The high accuracy obtained indicates that the selected features (N, P, K, temperature, humidity, pH, and rainfall) are highly predictive of optimal crop choice within the context of the dataset used. The attempt at feature engineering did not yield improvements, suggesting that the original features already capture the necessary information effectively for the models tested, or that more sophisticated feature engineering techniques might be required.

The performance of the Random Forest model, known for its ability to handle complex relationships and its resistance to overfitting (especially when tuned), makes it a suitable candidate for deployment in a real-world decision-support system. Such a system could empower farmers to make informed decisions, leading to optimized resource allocation, increased crop yields, and enhanced agricultural sustainability. The development of a React-based web application, as described in the project's README, serves as a practical example of how such a model can be deployed to provide accessible recommendations.

### 4.1. Limitations and Future Work

While the results are promising, certain limitations should be acknowledged. The dataset, though balanced, may not represent all possible agro-climatic zones or soil types. The model's performance is contingent on the quality and scope of the training data. Factors not included in this dataset, such as specific soil texture, pest and disease prevalence, irrigation facilities, and market demand for crops, could also influence optimal crop choice.

Future work could focus on:

- **Expanding the Dataset:** Incorporating data from diverse agro-ecological regions and including a wider array of features.
- **Exploring Advanced Models:** Investigating deep learning models or more complex ensemble techniques.

- **Real-time Data Integration:** Developing systems that can integrate real-time sensor data for dynamic recommendations.
- **Economic Analysis:** Incorporating economic factors like market prices and cultivation costs into the recommendation logic.
- **User Feedback and Iteration:** Deploying the tool (like the AgroAdvisor app) and gathering user feedback to iteratively improve the model and user experience.

## 5. Conclusion

This study successfully developed and evaluated a machine learning model for crop recommendation using essential environmental and soil parameters. The optimized Random Forest classifier demonstrated exceptional accuracy, precision, recall, and F1-score on unseen test data. The findings affirm the utility of machine learning as a powerful tool in precision agriculture, offering significant potential for enhancing decision-making processes for farmers. By providing data-driven crop recommendations, such systems can contribute to improved agricultural productivity, resource efficiency, and overall sustainability of farming practices. Further research and development, incorporating a broader range of data and features, can lead to even more robust and widely applicable crop recommendation systems.

## 6. References

[1] Placeholder, A. B., & Placeholder, C. D. (Year). *Title of a general agriculture book or paper*. Journal/Publisher.
[2] Placeholder, E. F. (Year). Soil-Climate-Crop Interactions. *Journal of Agronomy*.
[3] Placeholder, G. H., & Placeholder, I. J. (Year). Machine Learning in Precision Agriculture: A Review. *Computers and Electronics in Agriculture*.
[4] Scikit-learn: Machine Learning in Python, Pedregosa et al., JMLR 12, pp. 2825-2830, 2011.
[5] Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5-32.